

Early Predictor for Student Success Based on Behavioural and Demographical Indicators

```
In [234]: Import libraries

import zipfile
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns

Load data

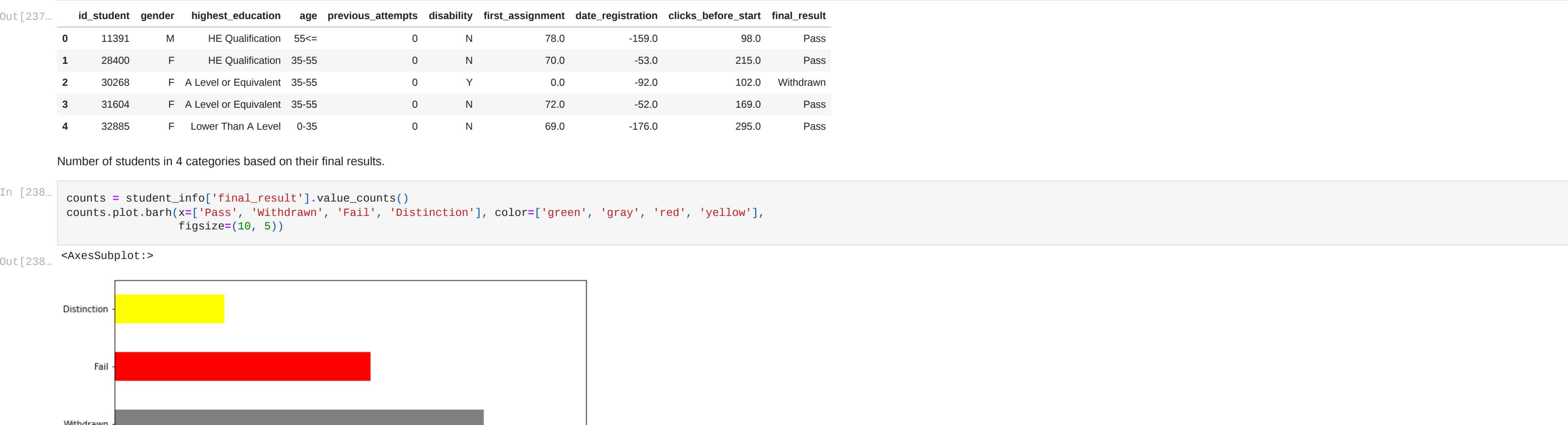
In [235]: zf = zipfile.ZipFile('./data.zip')
student_info = pd.read_csv(zf.open('studentInfo.csv'))
student_vle = pd.read_csv(zf.open('studentVle.csv'))
student_assessment = pd.read_csv(zf.open('studentAssessment.csv'))
student_registration = pd.read_csv(zf.open('studentRegistration.csv'))
assessments = pd.read_csv(zf.open('assessments.csv'))

Basic information about modules, first assessments and number of registered students

In [236]: # Get information about first day assessments
basic_info = assessments.filter(items=['code_module', 'code_presentation', 'date'])
basic_info = basic_info.groupby(['code_module', 'code_presentation']).size()
basic_info = pd.merge(basic_info, assessments, on=['code_module', 'code_presentation', 'date'],
                      how='inner').filter(items=['code_module', 'code_presentation', 'id_assessment', 'date'])

# Get total number of registered students for each module
registrations = student_registration.filter(items=['code_module', 'code_presentation', 'id_student'])
basic_info.insert(4, 'number of registered students', registrations['id_student'])
basic_info

Out[236]:
```



Converting all categorical variables into dichotomous variables.

```
In [239]: education_mapping = {
    'No formal quals': 0,
    'Lower Than A Level': 0,
    'A Level or Equivalent': 0,
    'HE Qualification': 1,
    'Post Graduate Qualification': 1
}

age_mapping = {
    '0-35': 0,
    '35-50': 1,
    '50+': 1
}

grade_mapping = {
    'Withdrawn': 0,
    'Fail': 1,
    'Pass': 2,
    'Distinction': 3
}

gender_mapping = {
    'F': 1,
    'M': 0
}

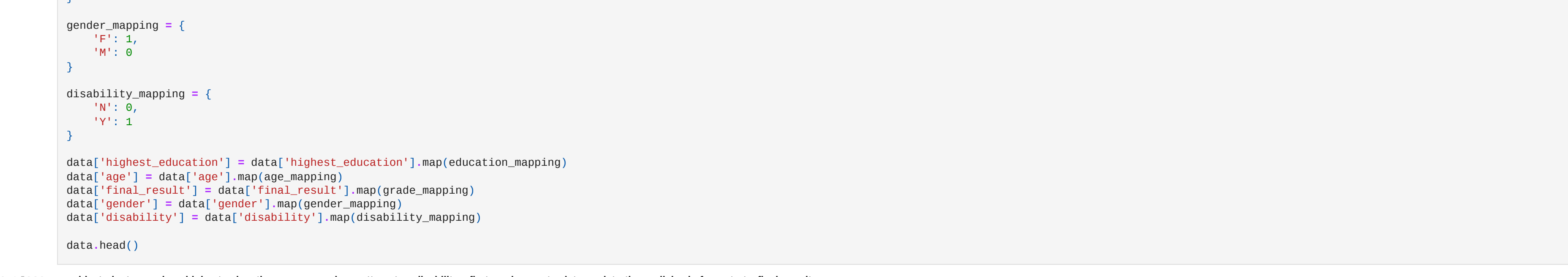
disability_mapping = {
    'N': 0,
    'Y': 1
}

data['highest_education'] = data['highest_education'].map(education_mapping)
data['age'] = data['age'].map(age_mapping)
data['final_result'] = data['final_result'].map(grade_mapping)
data['gender'] = data['gender'].map(gender_mapping)
data['disability'] = data['disability'].map(disability_mapping)
data.head()
```

Out[239]:

	id_student	gender	highest_education	age	previous_attempts	disability	first_assignment	date_registration	clicks_before_start	final_result
0	11391	0	1	1	0	0	78.0	-159.0	98.0	2
1	28400	1	1	1	0	0	70.0	-53.0	215.0	2
2	30268	1	0	1	0	1	0.0	-92.0	102.0	0
3	31004	1	0	1	0	0	72.0	-52.0	169.0	2
4	32885	1	0	0	0	0	69.0	-176.0	295.0	2

Pearson's correlation matrix for input variables.



Pearson chi-square test - testing whether the output was dependent upon the categorical input variables (Educational level, Age, Gender and Disability).

```
In [241]: from scipy.stats import chi2_contingency

test_df = pd.DataFrame(index=['statistic', 'p-value'])

for var in ['highest_education', 'age', 'gender', 'disability']:
    obs = pd.crosstab([var], data['final_result'])
    chi2, pval, _, _ = chi2_contingency(obs)
    test_df.insert(0, var, [chi2, pval])
test_df
```

Out[241]:

	disability	gender	age	highest_education
statistic	1.6456e+02	16.00200	2.6822e+02	3.0205e+02
pvalue	8.34139e-30	0.00083	1.45187e-64	1.08496e-02

Initialize lambdas that will be later used for one-vs-rest models.

```
In [242]: withdrawn_lambda = lambda x: 1 if x == 0 else -1
fail_lambda = lambda x: 1 if x == 1 else -1
pass_lambda = lambda x: 1 if x == 2 else -1
distinction_lambda = lambda x: 1 if x == 3 else -1
```

Generate training and testing sets.

```
In [243]: from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(data[['first_assignment', 'highest_education', 'age', 'gender', 'previous_attempts', 'disability', 'clicks_before_start', 'date_registration']],
                                                  data['final_result'], test_size=0.3, random_state=2)

y_train_withdrawn = y_train.map(withdrawn_lambda)
y_test_withdrawn = y_test.map(withdrawn_lambda)

y_train_fail = y_train.map(fail_lambda)
y_test_fail = y_test.map(fail_lambda)

y_train_pass = y_train.map(pass_lambda)
y_test_pass = y_test.map(pass_lambda)

y_train_distinction = y_train.map(distinction_lambda)
y_test_distinction = y_test.map(distinction_lambda)

y_trains = [y_train_withdrawn, y_train_fail, y_train_pass, y_train_distinction]
y_tests = [y_test_withdrawn, y_test_fail, y_test_pass, y_test_distinction]
```

Decision Tree

```
In [252]: from sklearn.tree import DecisionTreeClassifier

decision_trees = [DecisionTreeClassifier(class_weight='balanced', min_samples_split=10) for i in range(4)]

for index, tree in enumerate(decision_trees):
    tree.fit(X_train, y_trains[index])
```

Random Forest

```
In [253]: from sklearn.ensemble import RandomForestClassifier

random_forests = [RandomForestClassifier(n_jobs=-1, class_weight='balanced', bootstrap=False, min_samples_split=10) for i in range(4)]

for index, forest in enumerate(random_forests):
    forest.fit(X_train, y_trains[index])
```

Bayesian Additive Regression Trees (BART)

```
In [254]: from bartpy.sklearnmodel import SklearnModel

barts = [SklearnModel(n_burn=200, n_chains=2, n_samples=200, n_trees=50, alpha=0.9, beta=1.5) for i in range(4)]

for index, bart in enumerate(barts):
    bart.fit(X_train, y_trains[index])
```

Results

```
In [255]: from sklearn.metrics import accuracy_score, precision_score, f1_score, recall_score, confusion_matrix

results = []

for model, test in zip(decision_trees + random_forests, y_tests * 2):
    for score in [precision_score, recall_score, f1_score, accuracy_score]:
        results.append(score(y_test, model.predict(X_test)))

for model, test in zip(barts, y_tests):
    for score in [precision_score, recall_score, f1_score, accuracy_score]:
        results.append(score(y_test, model.predict(X_test)))

split = np.array_split(results, 3)

multi = pd.MultiIndex.from_product(
    [['Withdrawn', 'Fail', 'Pass', 'Distinction'], ['Precision', 'Recall', 'F1', 'Accuracy']],
    names=['final_result', 'Metric'])

results = pd.DataFrame(index=multi, columns=['Decision tree', 'Random Forest', 'BART'],
                      data=[decision_trees: split[0], 'Random Forest': split[1], 'BART': split[2]])
results
```

Out[255]:

		Decision tree	Random forest	BART
Withdrawn	Precision	0.49504	0.65797	0.774590
	Recall	0.664222	0.605959	0.560468
	F1	0.567587	0.629894	0.656377
	Accuracy	0.690837	0.782471	0.817895
Fail	Precision	0.254051	0.266097	0.333333
	Recall	0.416311	0.306620	0.028345
	F1	0.313544	0.297834	0.05642
	Accuracy	0.034952	0.060808	0.06868
Pass	Precision	0.934096	0.933432	0.937701
	Recall	0.971805	0.953208	0.991831
	F1	0.954147	0.946670	0.974351
	Accuracy	0.624054	0.642769	0.659542
Distinction	Precision	0.188446	0.257562	0.722322
	Recall	0.378905	0.316798	0.014656
	F1	0.251685	0.284125	0.028729
	Accuracy	0.795664	0.855195	0.910104



Vrsta istraživanja prema svrsi: primjenjeno istraživanje

- Želi se odrediti obrazac uspješna u učenju kako bi se moglo na vrijeme identificirati oni koji bi mogli imati poteškoće.

Vrsta istraživanja prema dubini: objašnjavno (explanatory) istraživanje

- Postavlja temeljne čimbenike i hipoteze o ključnim socio-ekonomskim i demografskim zračajcima koje utječu na uspjeh pojednca za daljnje, dublje analize.

Vrsta istraživanja prema vrsti korištenih podataka: kvantitativno istraživanje

- Podaci su kvantitativni te se koriste statističke metode i metode strojnog učenja kako bi se objasnio i predvidio uspjeh polaganja tečaja.

Vrsta istraživanja prema stupnju manipulacije varijabla: promatračko istraživanje (observacijsko, deskriptivno, neeksperimentalno)

- Nema intervencije istraživača, već se koriste podaci koji su skupljeni u prirodnom okruženju.

Vrsta istraživanja prema vrsti zaključaka: hipotetičko-deдуктивно istraživanje

- U ovom istraživanju promatra se stvarnost te pokušava hipoteza da socio-ekonomske i demografske karakteristike pojednca utječu na njihov uspjeh, zatim se pomoću modela strojnog učenja pokušavaju izolirati sve one karakteristike pomoću kojih je moguće zaključiti u predviđanje buduće uspješne i neuspješne.

Vrsta istraživanja prema vremenu provedbe: istraživanje iz perspektive ili sinkrono istraživanje

- Promatraju se određeni trenuci u istraživanju, konkretno početak tečaja i kraj tečaja. Nije naglasak na tome kako se pojedinci ponašaju tijekom tečaja, već što možemo reći o krajnjem rezultatu tečaja (prolaz, pad itd.) na temelju podataka koje imamo na početku, prije početka tečaja. Ne prate se kako se varijable mijenjaju tijekom tečaja jer većina njih (socio-ekonomske i demografske) se konstatirane o cijelo vrijeme.

Vrsta istraživanja prema izvoru informacija: sekundarno istraživanje

- Izvor informacija dolazi iz vanjskih izvora, OULAD (Open University Learning Analytics Dataset) te se želi identificirati obrasci i trendovi koji upućuju na drugačije ishode završetka tečaja.

Vrsta istraživanja prema načinu prikupljanje podataka: dokumentarno istraživanje

- Rad koristi vjerojatno bazu podataka koji su drugi istraživači prikupili, a za te istraživače možemo reći da su radili istraživanje jer su skupljali podatke sa izvora (tečaja).

Metode prikupljanja podataka: opservacijsko istraživanje

- Skupljeni podaci su dobiveni opservacijom, tj. promatraju stvari onako kako se događaju, bez ikakve intervencije. Radi se o sistematskom promatranju gdje se skupljaju kvantitativni objektivni podaci (poput prve ocjene ili broj klikova prije početka tečaja).

Istraživačka etika:

Na ovoj poveznici nađite se kratkai osvrt na način prikupljanja podataka. Najbitnije stvari su odabrane u nastavku.

- sudionici su pažljivo odabrani iz krstičnih i psiholoških odjeljaka
- Sudionici su anonimni i nisu bili kakve privatne informacije koje bi moglo identificirati studente)
- Sudionici su dobili prave privatne informacije about students and their modules. This includes the social security number, dates of birth and unique identifiers used at the OU for students. Module names have been replaced by semantic-free symbols and all temporal information has been expressed in relative terms with respect to the presentation start. In addition, all numeric identifiers (e.g., student_id, code_module, etc.) have been reassigned and completely randomised.

- povjerljivi podaci nisu uključeni u javni dataset